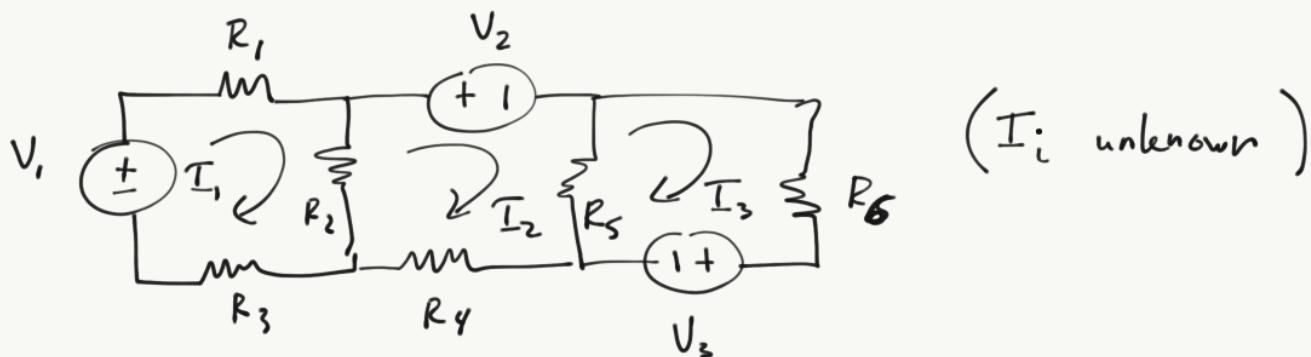


Linear Algebra

- Sets of linear equations play a key role in many physical systems → these can be represented as matrices
- For example analysis of circuits:



Applying Kirchoff's voltage law around each loop:

$$V_1 - I_1 R_1 - (I_1 - I_2) R_2 - I_1 R_3 = 0$$

$$-V_2 - (I_2 - I_3) R_5 - I_2 R_4 - (I_2 - I_1) R_2 = 0$$

$$- (I_3 - I_2) R_5 - I_3 R_6 - V_3 = 0$$

Rearranging, we find :

$$(-R_1 - R_2 - R_3) I_1 + R_2 I_2 + 0 = -V_1$$

$$+ R_2 I_1 (-R_5 - R_4 - R_2) I_2 + R_5 I_3 = V_2$$

$$+ 0 + R_5 I_2 (-R_5 - R_6) I_3 = V_3$$

- More generally a system of 3 eqns. & 3 unknowns :

$$A_{11}x_1 + A_{12}x_2 + A_{13}x_3 = b_1$$

$$A_{21}x_1 + A_{22}x_2 + A_{23}x_3 = b_2$$

$$A_{31}x_1 + A_{32}x_2 + A_{33}x_3 = b_3$$

- These can be represented in matrix form as :

$$\underline{\underline{A}} \underline{x} = \underline{b}$$

(note compactness!)

A_{ij} are elements of $\underline{\underline{A}}$
 x_j are elements of \underline{x}
 b_i are elements of \underline{b}

- Referring back to the system above & enforcing equivalence with matrix form, it is clear that matrix multiplication is defined by :

$$(\underline{\underline{A}} \underline{x})_i \equiv \sum_j A_{ij}x_j \rightarrow \text{more generally for two matrices}$$

$$(\underline{\underline{AB}})_{ij} \equiv \sum_k A_{ik}B_{kj}$$

matrix multiplication

- For multiplication to make sense matrices must be conformable, viz. of size: $(m \times p) \cdot \underbrace{(p \times n)}$

inner dimensions equal

$m \times n$

(3)

- Mathematically speaking, a system of equations could fall into one of four categories :
 - 1) no solution (inconsistent)
 - 2) trivial solution $\underline{x} = \underline{0}$
 - 3) infinite solutions conforming to some constraint
 - 4) unique solution
- For physical situations we expect 4) !!!
- Linear systems play a large role in numerical analysis - networks, fitting & estimation of parameters, interpolation, & ODEs/PDEs
- Special Matrices :

$$1) \text{ column vector } \underline{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \end{bmatrix}$$

$$2) \text{ row vector } \underline{y} = [y_1 \ y_2 \dots]$$

by notational default we assume \underline{x} is column &
 \underline{x}^T is row, unless otherwise stated ..

3) square matrix has size $n \times n$

$$4) \text{ diagonal matrix } \underline{\underline{D}} = \begin{bmatrix} d_{11} & & & \\ & d_{22} & & \\ & & \ddots & \\ & & & \ddots \end{bmatrix}$$

11

identity matrix $\underline{\underline{I}} = \begin{bmatrix} 1 & & & \\ & 1 & 0 & \\ & & \ddots & \\ 0 & & & 1 \end{bmatrix}$ (4)

upper triangular $\underline{\underline{U}} = \begin{bmatrix} u_{11} & u_{12} & u_{13} & \dots \\ & u_{22} & u_{23} & \\ 0 & & \ddots & \end{bmatrix}$

lower triangular $\underline{\underline{L}} = \begin{bmatrix} l_{11} & & & \\ l_{21} & l_{22} & & \\ l_{31} & l_{32} & l_{33} & \dots \end{bmatrix}$

banded $\underline{\underline{B}} = \begin{bmatrix} b_{11} & 0 & b_{13} & 0 \\ 0 & b_{22} & 0 & b_{24} \\ b_{31} & 0 & b_{33} & 0 \\ 0 & b_{42} & 0 & b_{44} \end{bmatrix}$

- Matrix operations all follow from linear equation properties

$$(\underline{\underline{A}} + \underline{\underline{B}})_{ij} = A_{ij} + B_{ij} ; (\underline{\underline{A}} - \underline{\underline{B}})_{ij} = A_{ij} - B_{ij}$$

$\underline{\underline{A}}, \underline{\underline{B}}$ must have same size!

$$(\alpha \underline{\underline{A}})_{ij} = \alpha A_{ij}$$

$$\underline{\underline{A}}(\underline{\underline{B}} + \underline{\underline{C}}) = \underline{\underline{AB}} + \underline{\underline{AC}} ; \quad \boxed{\underline{\underline{AB}} \neq \underline{\underline{BA}}}$$

(5)

- Matrix inverse $\underline{\underline{A}}^{-1}$ defined s.t. :

$$\underline{\underline{A}} \underline{\underline{A}}^{-1} = \underline{\underline{A}}^{-1} \underline{\underline{A}} = \underline{\underline{I}}$$

- Solutions to $\underline{\underline{A}} \underline{\underline{B}} = \underline{\underline{C}}$ for $\underline{\underline{A}}$:

$$\underline{\underline{A}} = \underline{\underline{B}}^{-1} \underline{\underline{C}}$$

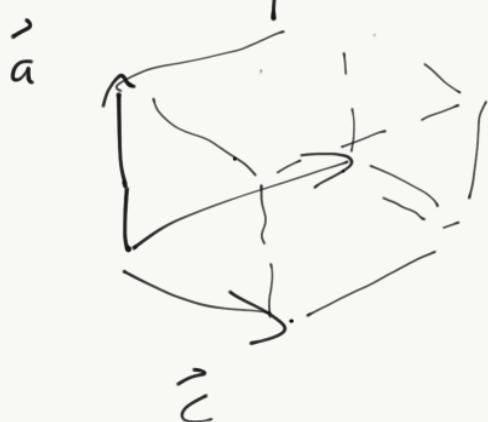
similarly $\underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{b}} \rightarrow \boxed{\underline{\underline{x}} = \underline{\underline{A}}^{-1} \underline{\underline{b}}}$

- Determinants of matrices can be thought of as analogous to concept of "volume" or magnitude

$$\underline{\underline{A}} = \begin{bmatrix} a_1 & a_2 & a_3 \\ b_1 & b_2 & b_3 \\ c_1 & c_2 & c_3 \end{bmatrix} \quad \vec{a} = (a_1, a_2, a_3) \\ \vec{b} = \dots \text{ etc.}$$

$\det(\underline{\underline{A}}) \equiv |\underline{\underline{A}}|$ is related to the triple product

$\vec{a} \cdot (\vec{b} \times \vec{c})$ which is volume spanned by $\vec{a}, \vec{b}, \vec{c}$



if $|\underline{\underline{A}}| = 0$

then \vec{a} lies
in plane of \vec{b}, \vec{c}

linearly dependent

- More generally we can define $|\underline{A}|$ by: ⑥

$$|\underline{A}| = \sum_{j=1}^n A_{ij} \text{ cof}(A_{ij}) = \sum_{j=1}^n A_{ij} (-1)^{i+j} M_{ij}$$

(i arb) \downarrow cofactor of A_{ij} \downarrow minor of A_{ij}

M_{ij} is the minor defined by evaluation det of augmented matrix, without row i , column j

- Note that $\det(\text{scalar}) = \text{scalar}$.

- Consider a 2×2 determinant

$$|\underline{A}| = \begin{vmatrix} A_{11} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = A_{11} (-1)^{1+1} M_{11} + A_{12} (-1)^{1+2} M_{12}$$

$$M_{11} \sim \begin{vmatrix} \cancel{A_{11}} & A_{12} \\ A_{21} & A_{22} \end{vmatrix} = \det(A_{22}) = A_{22}$$

$$M_{12} \sim \begin{vmatrix} A_{11} & \cancel{A_{12}} \\ A_{21} & \cancel{A_{22}} \end{vmatrix} = \det(A_{21}) = A_{21}$$

$$\boxed{|\underline{A}| = A_{11} A_{22} - A_{12} A_{21}} \quad (2 \times 2 \text{ det})$$

- Can also expand cofactors along column:

$$\boxed{|\underline{A}| = \sum_{i=1}^n A_{ij} \text{ cof}(A_{ij})} \quad (j \text{ fixed})$$

- 3×3 determinants proceed similarly : (7)

$$|A| = A_{11} (-1)^{1+1} M_{11} + A_{12} (-1)^{1+2} M_{12} + A_{13} (-1)^{1+3} M_{13}$$

$$M_{11} \sim \begin{vmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ A_{31} & A_{32} & A_{33} \end{vmatrix} = A_{22} A_{33} - A_{23} A_{32}$$

$$M_{12} = A_{21} A_{33} - A_{23} A_{31} ; \quad M_{13} = A_{21} A_{32} - A_{22} A_{31}$$

$$\begin{aligned} |A| = A_{11} (A_{22} A_{33} - A_{23} A_{32}) - A_{21} (A_{21} A_{33} - A_{23} A_{31}) \\ + A_{13} (A_{21} A_{32} - A_{22} A_{31}) \end{aligned}$$

- Evaluating an $n \times n$ determinant requires some adds
adds & $n!$ terms, each containing $(n-1)$
multiplications, so it scales like:

$$\mathcal{O}((n-1)n!) \quad \begin{matrix} \text{(cofactor expansion} \\ \text{of det)} \end{matrix}$$

Elimination Methods

(8)

- Cramer's rule is a direct method & useful for establishing a baseline approach for comparing efficiency of various methods:

$$\underline{\underline{A}} \underline{x} = \underline{b} ; \quad x_j = \frac{|\underline{\underline{A}}^j|}{|\underline{\underline{A}}|} ; \quad (\underline{\underline{A}}^j \text{ is } \underline{\underline{A}} \text{ with the } j^{\text{th}} \text{ column replaced by } \underline{b})$$

$$- \text{ e.g. } x_1 = \frac{\begin{vmatrix} b_1 & A_{12} \\ b_2 & A_{22} \end{vmatrix}}{A_{11}A_{22} - A_{12}A_{21}} ; \quad ; \quad \frac{\begin{vmatrix} A_{11} & b_1 \\ A_{21} & b_2 \end{vmatrix}}{A_{11}A_{22} - A_{12}A_{21}} = x_2$$

$$x_1 = \frac{b_1 A_{22} - A_{12} b_2}{A_{11}A_{22} - A_{12}A_{21}} ; \quad ; \quad x_2 = \frac{A_{11} b_2 - b_1 A_{21}}{A_{11}A_{22} - A_{12}A_{21}}$$

- Dealing with larger systems is problematic. For $\underline{\underline{A}} \sim n \times n$ we must evaluate $n+1$ determinants.
- Each determinant has $n!$ products (all containing n items multiplied) So $(n!)(n-1)$ multiplications per det.
- Total multiplications :

$$|\underline{\underline{M}}| = \mathcal{O}((n+1)(n-1)n!) \\ = \mathcal{O}((n-1)(n+1)!)$$

(9)

- Thus, Cramer's rule is too computationally expensive to practically use. (Slow IN MATLAB).
- Elimination methods address this shortcoming; these operate very much like a normal solution method by hand.
 - 1) perform algebraic substitutions until value of one variable can be found (fwd elimination)
 - 2) substitute this variable back into simplified eqns successively to find other vars.
(back substitution)
- Elementary row operations (allow elimination steps)
 - 1) Any row can be multiplied by scalar (scaling)
 - 2) Rows can be interchanged (pivoting)
 - 3) Rows may be linearly combined with other rows.
(elimination)
- None of these row operations change the solution to the system $\underline{A}\underline{x} = \underline{b}$
- Connection of row operations to matrix multiplication, e.g.:

$$1) \underline{E}_1 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & \alpha & 0 \\ 0 & 0 & 1 \end{bmatrix}; \underline{E}_1 \underline{A} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ \alpha A_{21} & \alpha A_{22} & \alpha A_{23} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}$$

$$2) \underline{E}_2 = \begin{bmatrix} 0 & 1 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 1 \end{bmatrix}; \underline{E}_2 \underline{A} = \begin{bmatrix} A_{21} & A_{22} & A_{23} \\ A_{11} & A_{12} & A_{13} \\ A_{31} & A_{32} & A_{33} \end{bmatrix}$$

$$\underline{\underline{E}}_3 = \begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ \alpha & 0 & 1 \end{bmatrix}; \quad \underline{\underline{E}}_3 \underline{\underline{A}} = \begin{bmatrix} A_{11} & A_{12} & A_{13} \\ A_{21} & A_{22} & A_{23} \\ \alpha A_{11} + A_{31} & \alpha A_{12} + A_{32} & \alpha A_{13} + A_{33} \end{bmatrix}$$

(10)

- Because of the connection of row ops to multiplication we can conceive as an inverse as a sequence of row ops.

$$\left. \begin{array}{l} \underline{\underline{A}}^{-1} \underline{\underline{A}} = \underline{\underline{I}} \\ \left(\prod_i \underline{\underline{E}}_i \right) \underline{\underline{A}} = \underline{\underline{I}} \end{array} \right\}$$

so its possible to
find an inverse
hence solution from
row operations!!!

- So if we find row ops that reduce $\underline{\underline{A}}$ to $\underline{\underline{I}}$ we can apply those ops to $\underline{\underline{I}}$ to get $\underline{\underline{A}}^{-1}$:

$$\underline{\underline{A}}^{-1} = \left(\prod_i \underline{\underline{E}}_i \right) \cdot \underbrace{\left(\prod_i \underline{\underline{E}}_i \right) \underline{\underline{I}}}_{\underline{\underline{I}}}$$

NOTE: Matrix multiplication is associative!

$$\underline{\underline{A}}(\underline{\underline{B}}\underline{\underline{C}}) = (\underline{\underline{A}}\underline{\underline{B}})\underline{\underline{C}}$$

$$\begin{array}{l}
 x_1 + 4x_2 + 2x_3 = 15 \\
 3x_1 + 2x_2 + x_3 = 10 \\
 2x_1 + x_2 + 3x_3 = 13
 \end{array} \rightarrow \left[\begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 3 & 2 & 1 & 10 \\
 2 & 1 & 3 & 13
 \end{array} \right] \quad (11)$$

elim. from row 1 :

$$\left[\begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 3 & 2 & 1 & 10 \\
 2 & 1 & 3 & 13
 \end{array} \right] \quad \begin{array}{l}
 R_2 - 3R_1 \\
 R_3 - 2R_1
 \end{array} \quad \begin{array}{c}
 \\ \\
 \cdots \cdots \cdots
 \end{array}$$

elim. from row R_2 :

$$\left[\begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 0 & -10 & -5 & -35 \\
 0 & -7 & -1 & -17
 \end{array} \right] \quad \begin{array}{l}
 R_3 - \frac{7}{10}R_2
 \end{array} \quad \begin{array}{c}
 \\ \\
 \begin{array}{l}
 -1 - \frac{7}{10}(-5) \\
 = \frac{5}{2}
 \end{array} \\
 -17 - \frac{7}{10}(-35)
 \end{array}$$

final state following elim :

$$\left[\begin{array}{ccc|c}
 1 & 4 & 2 & 15 \\
 0 & -10 & -5 & -35 \\
 0 & 0 & \frac{5}{2} & \frac{15}{2}
 \end{array} \right] \quad \begin{array}{c}
 \\ \\
 \begin{array}{l}
 -170 + 245 \\
 \hline
 10 \\
 = \frac{75}{10} = \frac{15}{2}
 \end{array}
 \end{array}$$

back substitution :

$$x_3 = \left(\frac{15}{2}\right) \left(\frac{2}{5}\right) = 3$$

$$x_2 = \left(-35 + 5(3)\right) \left(-\frac{1}{10}\right) = -20 \left(-\frac{1}{10}\right) = 2$$

$$x_1 = (15 - 2(3) - 4(2)) = 1$$

- Two modifications to standard elimination schemes are advisable : pivoting & slicing. (12)
- The pivot element for an elimination step is the element being used for elimination. e.g.

pivot $\left[\begin{array}{ccc|c} 3 & 2 & 1 & 10 \\ 2 & 1 & 3 & 13 \\ 1 & 4 & 2 & 15 \end{array} \right]$ $R_2 - \frac{2}{3}R_3$ $\xrightarrow{\text{elimination multiplier}}$
 $R_3 - \frac{1}{3}R_3$

- Interchanging rows (reordering eqns.) doesn't change solution - so one can switch pivots with impunity
- For reasons of numerical accuracy, it is always best to scale each row by largest element in that row & then switch to largest pivot.

$$\left[\begin{array}{ccc|c} 1 & 4 & 2 & 15 \\ 3 & 2 & 1 & 10 \\ 2 & 1 & 3 & 13 \end{array} \right] \quad \begin{aligned} \text{pivot: } & \frac{1}{15} \xrightarrow{\text{swap}} \\ \text{pivot: } & \frac{3}{10} \\ \text{pivot: } & \frac{2}{13} \end{aligned}$$

\downarrow

$$\left[\begin{array}{ccc|c} 3 & 2 & 1 & 10 \\ 1 & 4 & 2 & 15 \\ 2 & 1 & 3 & 13 \end{array} \right] \quad \begin{aligned} \text{do elimination, then reorder} \\ \text{again for largest scaled pivot} \end{aligned}$$

$\xrightarrow{\text{defines Gaussian elimination}}$

- Gauss-Jordan elimination is based on the fact that one may compute inverse as a sequence of row ops on $\underline{\underline{I}}$:

$$\underline{\underline{A}}^{-1} \underline{\underline{A}} = \underline{\underline{I}} \rightarrow \left(\prod_i \underline{\underline{E}}_i^{-1} \right) \underline{\underline{A}} = \underline{\underline{I}}$$

$$\therefore \left(\prod_i \underline{\underline{E}}_i^{-1} \right) \underline{\underline{I}} = \underline{\underline{A}}^{-1}$$

$$\left[\begin{array}{ccc|ccc} 1 & 4 & 2 & 1 & 0 & 0 \\ 3 & 2 & 1 & 0 & 1 & 0 \\ 2 & 1 & 3 & 0 & 0 & 1 \end{array} \right] = \left[\begin{array}{c|c} \underline{\underline{A}} & \underline{\underline{I}} \end{array} \right] \quad \begin{aligned} R_2 - 3R_1 \\ R_3 - 2R_1 \end{aligned}$$

- See book examples & homework.
- Generally elimination methods of all types are more efficient than solving a system via Cramer's Rule; which scaled like

$$\mathcal{O}(n(n-1)(n+1)!) \quad \underbrace{\text{Cramer's Rule}}$$

- E.g. Gauss elimination scales like: $\mathcal{O}\left(\frac{n^3}{3} + \frac{n}{3}\right)$

Determinants

(14)

- Determinants can be evaluated efficiently using elimination methods.
- Consider an upper triangular matrix resulting from elimination:

$$\underline{U} = \begin{bmatrix} U_{11} & U_{12} & U_{13} & \dots \\ 0 & U_{22} & U_{23} & \dots \\ 0 & 0 & U_{33} & \dots \\ 0 & 0 & 0 & \dots \end{bmatrix}$$

$\det(\underline{U}) = (\text{expand along } 1^{\text{st}} \text{ column})$

$$= U_{11} (-1)^{1+1} \left| \underline{M}_{11} \right| = U_{11} (-1)^{1+1} \left[U_{22} (-1)^{2+1} \left| \underline{M}_{22} \right| \right]$$

$$= \dots = \prod_i U_{ii} \quad (\text{expand along } 1^{\text{st}} \text{ column of } \underline{M}_{ii})$$

- So evaluating dets of triangular matrices is trivial
- Performing simple elimination does not change the determinant, thus we can do fwd elimination then evaluate det from result.
- This is more efficient than the $(n-1)n!$ operations implied by cofactor expansions.
- Pivoting changes value of determinants in predictable ways:

- 1) scaling row by a const also scales det by const
- 2) swapping rows flips sign of det.

LU factorization

(15)

- One can factor a matrix in infinite ways; one useful way is to factor into upper & lower triangular matrices:

$$\underline{\underline{A}} = \underline{\underline{L}} \underline{\underline{U}}$$

- This can be used to solve $\underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{b}}$ as follows:

$$\underline{\underline{L}} \underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{b}} ; \quad \underline{\underline{L}}^{-1} \underline{\underline{L}} \underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{L}}^{-1} \underline{\underline{b}}$$

$$\underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{L}}^{-1} \underline{\underline{b}}$$

- Defining $\underline{\underline{b}}' = \underline{\underline{L}}^{-1} \underline{\underline{b}}$ we have the following eqns:

$$(1) \quad \underline{\underline{L}} \underline{\underline{b}}' = \underline{\underline{b}} \quad (2) \quad \underline{\underline{U}} \underline{\underline{x}} = \underline{\underline{b}}'$$

- Thus we can view solving $\underline{\underline{A}} \underline{\underline{x}} = \underline{\underline{b}}$ as solving (1), (2) above
- This is advantageous since triangular systems can be solved easily using forward/backward substitution. Also for systems with multiple RHS, $\underline{\underline{L}}, \underline{\underline{U}}$ need only be computed once!
- An example with multiple RHS would be finding an inverse via LU factorization:

find $\underline{\underline{A}}^{-1}$ s.t. $\underline{\underline{A}} (\underline{\underline{A}}^{-1}) = \underline{\underline{I}}$ can be viewed as

76

$$\left. \begin{array}{l} \underline{\underline{A}} \underline{\underline{x}}_1 = \underline{\underline{b}}_1 \\ \underline{\underline{A}} \underline{\underline{x}}_2 = \underline{\underline{b}}_2 \\ \underline{\underline{A}} \underline{\underline{x}}_3 = \underline{\underline{b}}_3 \end{array} \right\} \quad \underline{\underline{A}}^{-1} = \begin{bmatrix} \underline{\underline{x}}_1 & \underline{\underline{x}}_2 & \underline{\underline{x}}_3 \end{bmatrix}$$
$$\underline{\underline{b}}_1 = \underline{\underline{I}}_{i1} \quad (\text{first column of } \underline{\underline{I}})$$
$$\underline{\underline{b}}_2 = \underline{\underline{I}}_{i2} \quad (2^{\text{nd}} \quad " \quad " \quad ")$$
$$\underline{\underline{b}}_3 = \underline{\underline{I}}_{i3} \quad (3^{\text{rd}} \quad " \quad " \quad ")$$

Tri diagonal Systems

(17)

- Consider the system :

$$\underline{\underline{T}} \underline{x} = \underline{b};$$

$$\underline{\underline{T}} = \begin{bmatrix} T_{11} & T_{12} & & & \\ T_{21} & T_{22} & T_{23} & & \\ & T_{23} & T_{33} & T_{34} & \\ & & T_{34} & T_{44} & T_{45} \\ & & & \ddots & \\ & & & & T_{i,i-1} & T_{ii} & T_{i,i+1} \\ & & & & & \ddots & \\ & & & & & & T_{n,n-1} & T_{nn} \end{bmatrix}$$

- Elimination applied to this system is particularly simple. E.g. only one element in row 2, to be eliminated using row 1. (column 1)

$$R_2 \rightarrow R_2 - \frac{T_{21}}{T_{11}} R_1; \text{ viz. } \underline{\underline{T}} \underline{x} = \underline{b} \quad \underbrace{\qquad}_{j \in \{2,3\}}$$

$$T_{2j} = T_{2j} - \frac{T_{21}}{T_{11}} T_{1j}$$

- Now full elimination can proceed downward to subsequent rows ...

- We must also alter RHS from elim step : (18)

$$b_2 = b_2 - \frac{T_{21}}{T_{11}} b_1$$

- These two steps can be encoded generally as ($j \leq n$)

$$1) \quad T_{i,j} = T_{i,j} - \left(\frac{T_{i,i-1}}{T_{i-1,i-1}} \right) T_{i-1,j}$$

$$(i \in \{1, n\}; j \in \{i-1, i+1\})$$

there are two of these to perform :

$$T_{i,i} = T_{i,i} - \left(\frac{T_{i,i-1}}{T_{i-1,i-1}} \right) T_{i-1,i}$$

$$T_{i,i+1} = T_{i,i+1} - \left(\frac{T_{i,i-1}}{T_{i-1,i-1}} \right) T_{i-1,i+1}$$

- 2) Following by RHS :

$$b_i = b_i - \left(\frac{T_{i,i-1}}{T_{i-1,i-1}} \right) b_{i-1}$$

- Following these two steps \underline{T} has been triangulated

$\underline{T} \rightarrow \underline{T}'$; we can now perform backsubstitution...

$$x_n = \frac{b_n}{T_{nn}}; x_i = \frac{1}{T_{ii}} (b_i - T_{i,i+1} x_{i+1})$$

Condition Number

19

- Elimination methods can be sensitive to round-off error which propagates thru subsequent steps ...
- One way to measure this sensitivity is thru the matrix condition number which is computed using norms
- A norm is a measure of magnitude & satisfies the constraints:

$$a) \|\underline{A}\| > 0$$

$$b) \|\underline{A}\| = 0 \text{ iff } \underline{A} = \underline{0}$$

$$c) \|\alpha \underline{A}\| = |\alpha| \|\underline{A}\|$$

$$d) \|\underline{A} + \underline{B}\| = \|\underline{A}\| + \|\underline{B}\|$$

$$e) \|\underline{A}\underline{B}\| \leq \|\underline{A}\| \|\underline{B}\| \quad (\underbrace{\text{Schwarz}}_{\text{inequality}})$$

- Commonly used vector norms:

$$\|\underline{x}\|_1 = \sum_i |x_i| \quad L^1 - \text{norm}$$

$$\|\underline{x}\|_2 = \sqrt{\sum_i x_i^2} \quad L^2 - \text{norm} \quad (\text{Euclidean})$$

$$\|\underline{x}\|_\infty = \max(|x_i|) \quad L^\infty - \text{norm}$$

- Matrix norms can be similarly defined : (20)

$$\|\underline{A}\|_1 = \max_j \left\{ \sum_i |A_{ij}| \right\} \quad (\text{max col. sum})$$

$$\|\underline{A}\|_\infty = \max_i \left\{ \sum_j A_{ij} \right\} \quad (\text{max row sum})$$

$$\|\underline{A}\|_2 = \min(\lambda_{\cdot\cdot}) \quad / \begin{matrix} \text{min eigenvalue,} \\ \text{spectral norm} \end{matrix}$$

$$\|\underline{A}\|_e = \sqrt{\sum_{i,j} A_{ij}^2} \quad (\text{Euclidean norm})$$

- The condition number measures sensitivity to small changes (e.g. round-off) in system.
- Considering the system :

$$\underline{A} \underline{x} = \underline{b}$$

we know $\|\underline{b}\| \leq \|\underline{A}\| \|\underline{x}\|$

- Let us examine system perturbed by $\underline{\delta b}$

$$\underline{A}(\underline{x} + \underline{\delta x}) = \underline{b} + \underline{\delta b} = \underline{A}\underline{x} + \underline{A}\underline{\delta x} = \underline{b} + \underline{\delta b}$$

- Subtract original equation :

$$\underline{A} \underline{\delta x} = \underline{\delta b}$$

$$\Rightarrow \underline{\delta x} = \underline{A}^{-1} \underline{\delta b}$$

thus $\|\underline{\delta x}\| \leq \|\underline{A}^{-1}\| \|\underline{\delta b}\|$

- Multiplying the two Schwarz-inequalities : (21)

$$\|\underline{b}\| \|\underline{\delta x}\| \leq \|\underline{A}\| \|\underline{x}\| \|\underline{A}^{-1}\| \|\underline{\delta b}\|$$

$$\|\underline{b}\| \|\underline{\delta x}\| \leq \|\underline{x}\| \|\underline{\delta b}\| (\|\underline{A}\| \|\underline{A}^{-1}\|)$$

$$\frac{\|\underline{\delta x}\|}{\|\underline{x}\|} \leq \frac{\|\underline{\delta b}\|}{\|\underline{b}\|} C(\underline{A})$$

$$C(\underline{A}) \equiv \|\underline{A}\| \|\underline{A}^{-1}\| \quad (\text{condition number})$$

- A large condition number \Rightarrow modest $\|\underline{\delta b}\|/\|\underline{b}\|$ can lead to large $\|\underline{\delta x}\|/\|\underline{x}\|$ (BAD!!)
- $C(\underline{A}) \sim 1$ is best.

Iterative Methods

(22)

- One drawback of elimination methods is the accumulation of error through successive steps.
- Iterative methods can address this issue, but require a matrix to be diagonally dominant, i.e.

$$\boxed{\exists i \text{ s.t. } |A_{ii}| \geq \sum_{j, j \neq i} |A_{ij}|}$$

- Many Systems will not meet this requirement...
- Iterative methods start with an initial guess and then refine that guess repeatedly until convergence is achieved.
- Consider: $\underline{A}\underline{x} = \underline{b}$, for some guess $\underline{x}^{(0)}$

the residual (error in a sense) is $\underline{b} - \underline{A}\underline{x}^{(0)}$

$$\boxed{R_i^{(k)} = b_i - \sum_j A_{ij}x_j^{(k)}} \quad \begin{matrix} \text{residual on } i^{\text{th}} \\ \text{var } \in \text{ iteration} \\ k \end{matrix}$$

- When the residual is small the iterated solution is near the true solution.
- When the solution $\underline{x}^{(k)}$ does not change appreciably between iterations, it has converged.

- Jacobi iteration involves a correction to each variable based on residual : (23)

$$x_i^{(k+1)} = x_i^{(k)} + \frac{r_i^{(k)}}{A_{ii}}$$

- The basis of this idea is that, given some guess for all vars. (sans i) we can directly compute the variable i :

$$\begin{aligned} A_{ii}x_i^{(k+1)} &\approx b_i - \sum_{j \neq i} A_{ij}x_j^{(k)} \\ &= A_{ii}x_i^{(k)} + b_i - \sum_j A_{ij}x_j^{(k)} \\ \therefore x_i^{(k+1)} &= x_i^{(k)} = \frac{1}{A_{ii}} \left(b_i - \sum_j A_{ij}x_j^{(k)} \right) \end{aligned}$$

- letting $\underline{\Delta x} = \underline{x}^{(k+1)} - \underline{x}^{(k)}$ we can define convergence as : $|\underline{\Delta x}| \leq \varepsilon$ for some iteration step
- Relative error can also be used for convergence but becomes problematic if a solution is near zero.

- Gauss-Seidel iteration uses the partially updated solution to determine $R_i^{(k)}$:

$$R_i^{(k)} = b_i - \sum_{j < i} A_{ij} x_j^{(k+1)} - \sum_{j \geq i} A_{ij} x_j^{(k)}$$

$$x_i^{(k+1)} = x_i^{(k)} + \frac{R_i^{(k)}}{A_{ii}} \quad (\text{same})$$

- Convergence is almost always faster than Jacobi iteration & it's "simpler" to code.
- Successive over-relaxation (SOR) is a further refinement to speed convergence.

$$x_i^{(k+1)} = x_i^{(k)} + \omega \frac{R_i^{(k)}}{A_{ii}} \quad (\omega < 2)$$

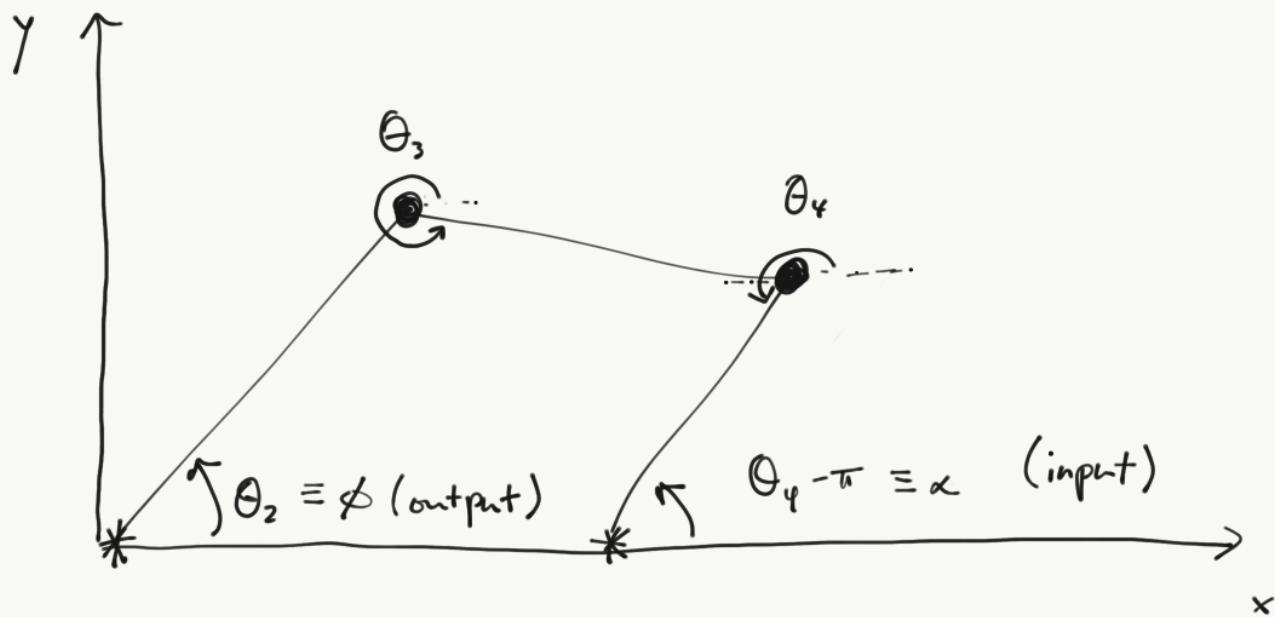
- The relaxation parameter (ω) must be chosen essentially by hand & some choices will make convergence slower...

Nonlinear Equations

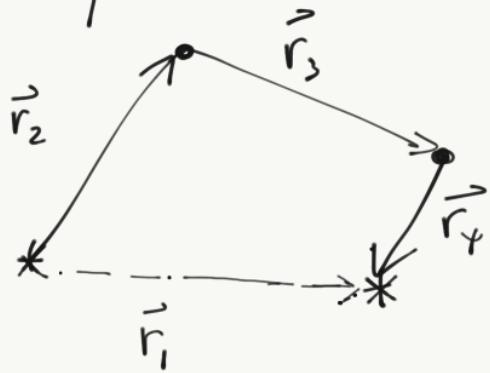
- Solutions of nonlinear equations are necessary in a wide range of applications.
- The canonical nonlinear problem is that of finding the roots of a function. I.e find

$$\left\{ \begin{array}{l} \{x_i\} \\ \text{s.t. } f(x_i) = 0 \end{array} \right.$$

- The book gives the example of a four bar linkage :



Vectorially we have :



$$\vec{r}_1 = \vec{r}_2 + \vec{r}_3 + \vec{r}_4$$

$$(x) \quad r_1 = r_2 \cos \theta_2 + r_3 \cos \theta_3 + r_4 \cos \theta_4$$

$$(y) \quad 0 = r_2 \sin \theta_2 + r_3 \sin \theta_3 + r_4 \sin \theta_4$$

- Nonlinear equations can be extremely difficult to ② solve by hand, e.g.:

$$x - e^{-x} = 0 \quad (\text{transcendental})$$

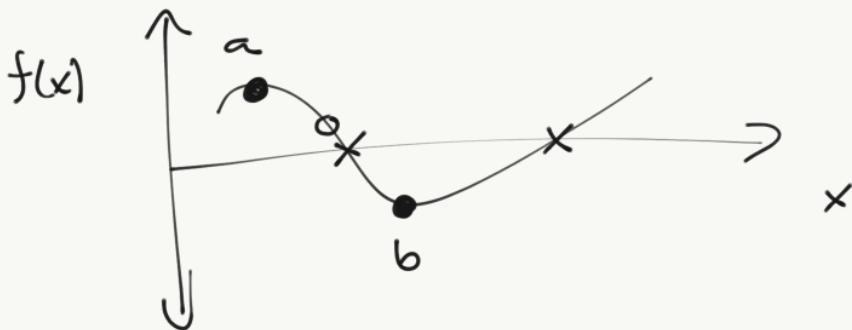
- Often you just have to plot things or start guessing ...

Root finding

(3)

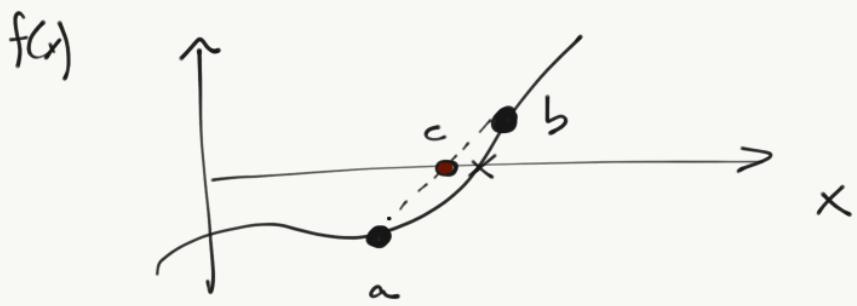
- Root finding procedures (finding zeros of $f(x)$) generally share several basic steps:
 - 1) bounding (restriction of solution to some interval, $a \leq x_i \leq b$)
 - 2) refinement (update solution to reduce error)
- Bounding
 - a) graphing or search
 - b) prior knowledge
 - c) solution of approximate model
 - d) prior solution in sequence of iterations.
- Refinement
 - a) guessing
 - b) closed domain methods
 - c) open domain "
- See figure 3.6 & 3.2.4
- Root finding algorithm features
 - a) max # of iterations
 - b) error checking on whether $f'(x) \rightarrow 0$
 - c) convergence test (x_i not changing or $f(x_i)$ small)
 - d) checking that the root actually works.

- Interval halving is a straightforward way to implement a search method for root-finding
- If the value of the function at two different points changes sign, & that function is continuous then the root lies between those points:



- Just guessing one could estimate that the root is halfway in between
$$x_i \sim \frac{a+b}{2}$$
- We can then choose the half interval (a, c) or (c, b) that contains the root (fn. changes sign.) & iterate until $f(c) = 0$ to some desired tolerance.
- False position is a refinement of interval halving.
- Rather than taking c @ the midpoint of (a, b) , we approximate a line connecting $f(a)$ & $f(b)$ & find where that line crosses zero to estimate c .

(5)



(value of fn.
at left
↑ "boundary")

→ The line is defined by $y = m(x-c) + y_a$

$$m \approx \frac{f(b) - f(a)}{b-a}; \quad y(x) = \frac{f(b) - f(a)}{b-a}(x-a) - f(a)$$

c: $0 = \frac{f(b) - f(a)}{b-a} (c-a) + f(a)$

↑

(position @ which $y(x) = 0$) \therefore

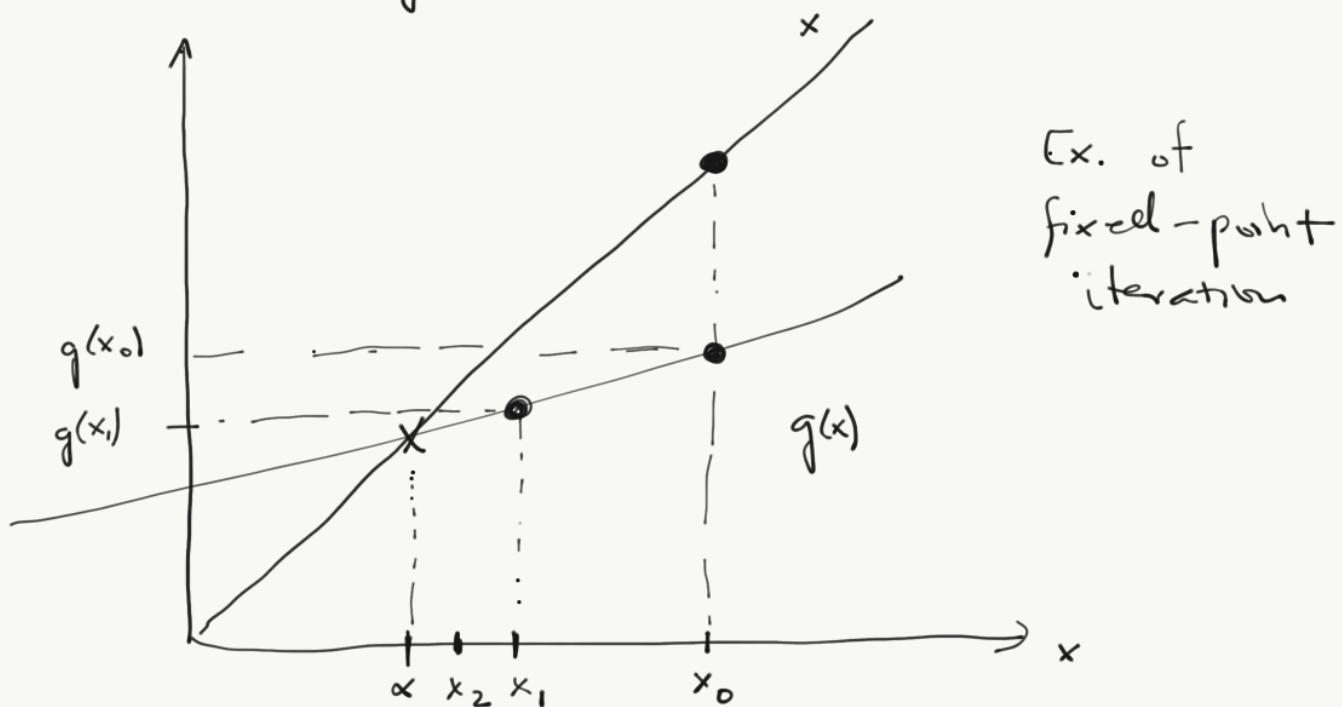
$$c = a - \frac{f(a)}{\frac{(b-a)}{f(b) - f(a)}}$$

Now as before we choose a new interval containing the root & iterate to convergence.

Fixed-point iteration

(6)

- Open domain methods do not require a root to be "bracketed"
 - a) fixed-point iteration
 - b) Newton's method
 - c) secant method
- Fixed point iteration:
 - 1) recast $f(x) = 0$ into $x = g(x)$
 - 2) guess x_0 (initial root estimate)
 - 3) iterate to convergence :
$$x_{i+1} = g(x_i)$$
- This is similar in concept to Jacobi iterations for linear equations - guess solution & use it to solve remaining vars.



- Iterations (of any method) continue until some convergence is achieved. ⑦

$$|x_{i+1} - x_i| \leq \varepsilon \quad \text{or} \quad |f(x_{i+1})| \leq \varepsilon'$$

- Suppose the true root (unknown) is denoted by α ; the error is then $x - \alpha$.

$$x_{i+1} = g(x_i) \quad \therefore \quad x_{i+1} - \alpha = g(x_i) - \alpha$$

- At the solution $\alpha = g(\alpha)$:

$$x_{i+1} - \alpha = g(x_i) - g(\alpha)$$

thus: $e_{i+1} = g(x_i) - g(\alpha)$

- let us expand $g(\alpha)$ in a Taylor Series: (about x_i)

$$g(\alpha) = g(x_i) + g'(x_i)(\alpha - x_i) + \dots$$

then: $e_{i+1} \approx g(x_i) - \{g(x_i) + g'(x_i)(\alpha - x_i)\}$

$$\approx -g'(x_i) \underbrace{(\alpha - x_i)}_{e_i} \quad \begin{array}{l} \text{(linear convergence} \\ \text{e_{i+1} is linear fn.} \\ \text{of e_i}) \end{array}$$

$$\Rightarrow \boxed{\begin{aligned} e_{i+1} &= \overline{g'(x_i)} e_i \Rightarrow |e_{i+1}| = |g'(x_i)| |e_i| \\ |e_{i+1}| / |e_i| &= |g'(x_i)| \text{ (must be } < 1) \end{aligned}}$$

(8)

Newton's Method

- Converges quadratically, but requires derivative eval.

$$f'(x_i) \approx \frac{f(x_{i+1}) - f(x_i)}{x_{i+1} - x_i} \quad (\text{truncated Taylor Series})$$

- Solving for x_{i+1} , we find:

$$f'(x_i)(x_{i+1} - x_i) = f_{i+1} - f_i$$

$$f'_i x_{i+1} = f_{i+1} - f_i + f'_i x_i$$

$$x_{i+1} = \frac{f_{i+1} - f_i}{f'_i} + x_i$$

- The purpose of an iteration like this is to drive the root to zero, viz. $f_{i+1} = 0$

$$\therefore \boxed{x_{i+1} = x_i - \frac{f(x_i)}{f'(x_i)}} \quad \text{Newton method!}$$

- Convergence can be quantified by noting similarity to fixed-point iteration:

$$x_{i+1} = g(x_i) \quad w/ \quad g(x_i) = x_i - \frac{f(x_i)}{f'(x_i)}$$

- For fixed point $|g'(x_i)| < 1$ converges.

$$\frac{d}{dx} \left(x - \frac{f(x)}{f'(x)} \right) = 1 - \frac{f'f' - f''f}{(f')^2} \quad (9)$$

$$g'(x) = \frac{(f')^2 - (ff')^2 + f''f}{(f')^2} = \frac{f''f}{(f')^2}$$

- Evaluating @ root $x = \alpha$ \rightarrow (goes to zero @ $x = \alpha$)

$$g'(\alpha) = \frac{f''(\alpha)f(\alpha)}{(f')^2} = 0 \quad (\text{convergent!})$$

- Similar to before $e = x - \alpha$

$$\begin{aligned} e_{i+1} &= x_{i+1} - \alpha = x_i - \frac{f_i}{f'_i} - \alpha \\ &= x_i - \alpha - \frac{f_i}{f'_i} = e_i - \frac{f_i}{f'_i} \end{aligned}$$

- Expanding about x_i :

$$f(x) = f(x_i) + f'(x_i)(\alpha - x_i) + \frac{1}{2} f''(x_i)(\alpha - x_i)^2 + \dots$$

- At root $f(x) = 0 \Rightarrow$

$$f(x_i) = f'(x_i)(x_i - \alpha) - \frac{1}{2} f''(x_i)(x_i - \alpha)^2$$

$$= f'(x_i)e_i - \frac{1}{2} f''(x_i)e_i^2$$

$$\therefore e_{i+1} = e_i - \frac{f'(x_i)e_i - \frac{1}{2} f''(x_i)e_i^2}{f'(x_i)} = \frac{\frac{1}{2} \frac{f''(x_i)}{f'(x_i)} e_i^2}{\underline{f'(x_i)}}$$

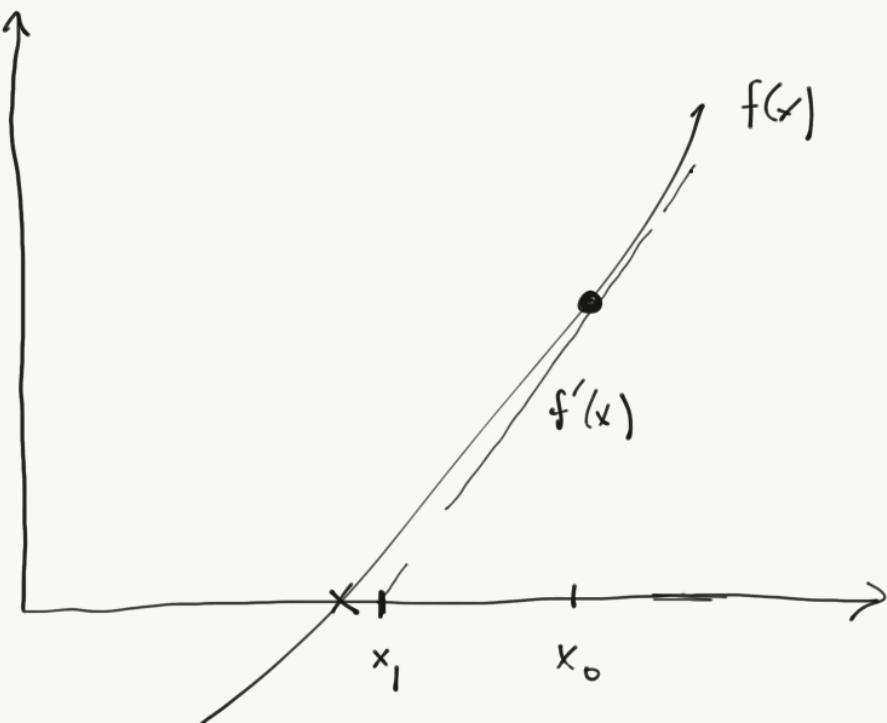
(10)

- In the case that derivatives are too expensive to evaluate they can be approximated:

$$\tilde{f}'(x_i) \approx \frac{f(x_{i+1}) - f(x_i)}{(x_{i+1} - x_i)}$$

$$x_{i+1} = x_i - \frac{f(x_i)}{\tilde{f}'(x_i)} \quad (\text{approx. Newton's method})$$

Ex. of
Newton's
method



Secant Method

- Bears a lot of similarity to approximate Newton's method
- Given two initial approximations of the root (x_0, x_1), we take :

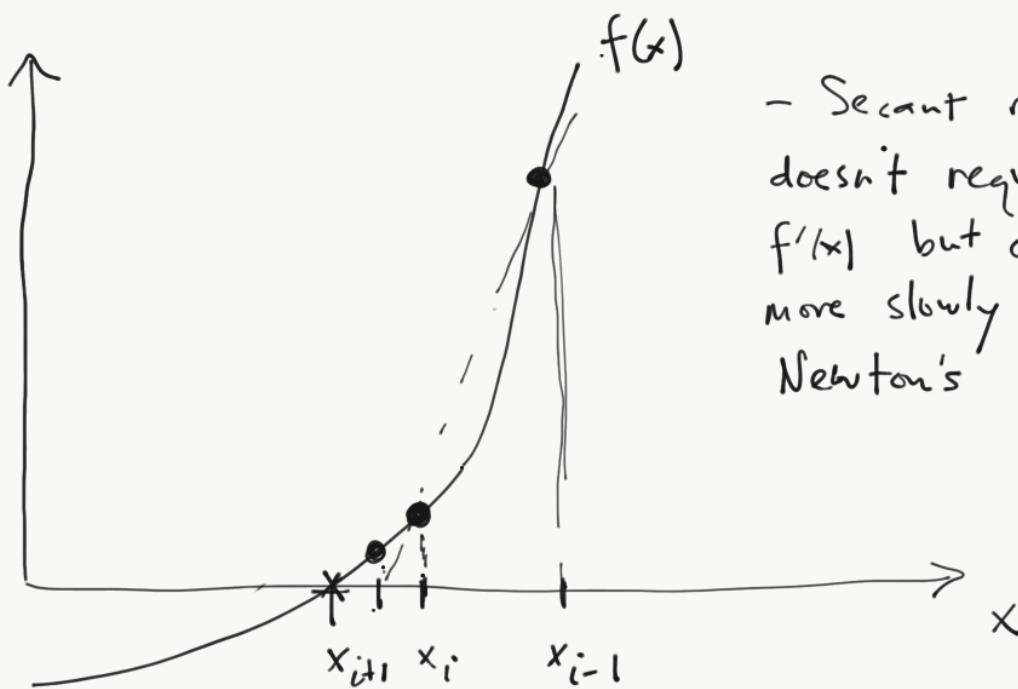
$$g'(x_i) \approx f'(x_i) = \frac{f(x_i) - f(x_{i-1})}{x_i - x_{i-1}}$$

- This derivative is used w/ Newton-style iteration:

$$x_{i+1} = x_i - \frac{f(x_i)}{g'(x_i)}$$

- This is equivalent to finding root of secant line as next approx. to root of $f(x)$.

(Ex. of
secant
method)



- Secant method doesn't require $f'(x)$ but converges more slowly than Newton's

Muller's Method

(12)

- Muller's method takes 3 initial points & approximates next root iteratively from a quadratic form fitting these.

- Given x_i, x_{i-1}, x_{i-2} we find x_{i+1} from :

$$g(x) = a(x-x_i)^2 + b(x-x_i) + c$$

- With a, b, c determined from enforcing $g(x)$ passes thru $f @ x_i, x_{i-1}, x_{i-2}$

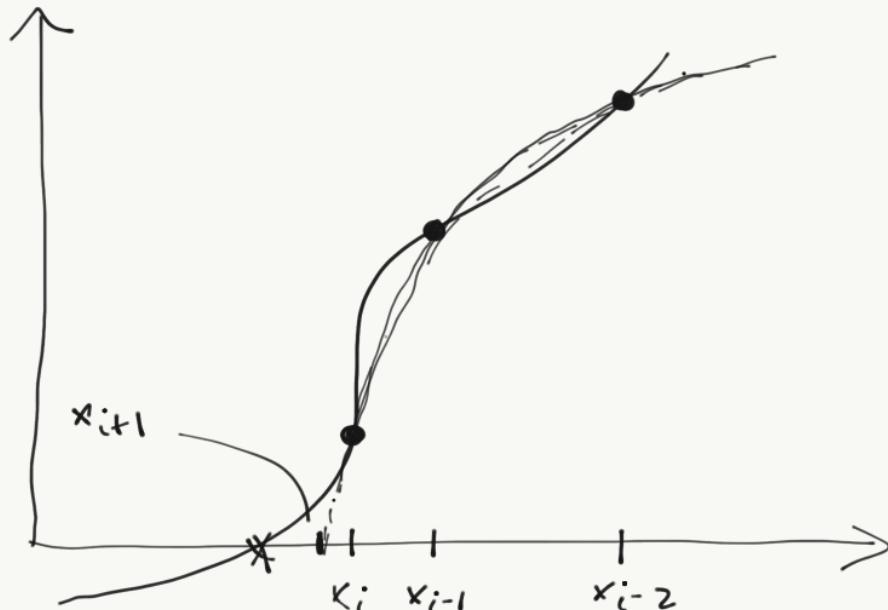
$$f_i = c$$

$$f_{i-1} = a(x_{i-1}-x_i)^2 + b(x_{i-1}-x_i) + c$$

$$f_{i-2} = a(x_{i-2}-x_i)^2 + b(x_{i-2}-x_i) + c$$

- a, b solved using latter two eqns. (see book)

Ex: Muller's method.



(13)

- Once a, b, c known, we find root of quadratic to estimate next value of root of $f(x)$:

$$x_{i+1} = x_i - \frac{2c}{b \pm \sqrt{b^2 - 4ac}}$$

- We choose root having $(\pm) \rightarrow \text{sign}(b)$ to prevent overly large changes $\Delta x = x_{i+1} - x_i$

Polynomial roots

(contains material
from 4.2)

(14)

- Polynomials are of the form :

$$P_n(x) = a_0 + a_1x + a_2x^2 + \dots + a_nx^n$$

- Fundamental thm of algebra: n^{th} degree polynomial has n roots.
- For linear & quadratics eqns can solve roots directly:

$$P_1(x) = ax + b \Rightarrow x = -\frac{b}{a}$$

$$P_2(x) = ax^2 + bx + c \Rightarrow x = \frac{-b \pm \sqrt{b^2 - 4ac}}{2a}$$

if roots have Im part they occur in complex conjugate pairs.

- Descartes rule of signs: number of positive roots of

$P_n(x)$ is number of sign changes in nonzero coeffs or smaller by an even integer. Number of negative roots found by considering $P_n(-x)$ in the same manner.

- Root-finding tends to be ill-conditioned in that small changes to P_n coeffs can greatly change roots.
- Polynomial evaluations require lots of multiplications due to x^2, x^3, x^4 and so on.

(15)

- Nested multiplication can reduce number of operations needed to evaluate polynomial e.g.

$$P_4(x) = a_0 + a_1x + a_2x^2 + a_3x^3 + a_4x^4$$

$$= a_0 + x(a_1 + x(a_2 + x(a_3 + x a_4)))$$

- More generally, nested multiplication goes like :

$$P_n(x) = a_0 + x(a_1 + x(a_2 + \dots + x(a_{n-1} + x a_n) \dots))$$

which can be evaluated thru the sequence :

$$b_n = a_n$$

$$b_i = a_i + x b_{i+1} \quad (i = n-1, 0)$$

$$\Rightarrow \underline{\underline{b_0}} = P_n(x)$$

- Polynomial division is defined by :

$$P_n(x) = (x-N) Q_{n-1}(x) + R$$

\downarrow \downarrow \swarrow
 divisor quotient remainder
 (degree $n-1$)

$$\underline{\underline{P_n(N)}} = R$$

note $P_n(N) = 0 \Rightarrow (x-N)$ is a factor of $P_n(x)$
 $\Rightarrow N$ is a root

- From poly-division we have the following derivative constraint : (16)

$$P_n'(x) = Q_{n-1}(x) + (x-N) Q'_{n-1}(x)$$

- Therefore, at the root N :

$$P_n'(N) = Q_{n-1}(N)$$

$$Q_{n-1} = \sum_{i=1}^{n-1} b_i x^{i-1}$$

$$P_n(x) = \sum_{i=0}^n a_i x^i$$

- Given a number, N , we can perform division as follows :

$$P_n(x) = (x-N) Q_{n-1}(x) + R$$

$$\sum_{i=0}^n a_i x^i = (x-N) \sum_{i=1}^{n-1} b_i x^{i-1} + R$$

$$\sum_{i=0}^n a_i x^i = \sum_{i=1}^{n-1} b_i x^i - \left(\sum_{i=1}^{n-1} b_i x^{i-1} \right) N + R$$

- Polynomial representations are unique so we can equate coefficients of the powers :

$$i' < i-1$$

$$i = i'+1$$

$$i = 1 \Rightarrow i' = 0$$

$$i = n-1 \Rightarrow i' = n-2$$

$$\sum_{i=0}^n a_i x^i = \sum_{i=1}^{n-1} b_i x^i - N \left(\sum_{i=0}^{n-2} b_{i+1} x^i \right) + R$$

$$a_i = b_i - N b_{i+1} \Rightarrow b_i = a_i + N b_{i+1}$$

- First & last terms are special cases:

(17)

$$b_n = a_n$$

$$\overline{b_0 = R = a_0 + Nb_1}$$

- So synthetic division can be summarized as:

$$\boxed{\begin{aligned} b_n &= a \\ b_i &= a_i + Nb_{i+1} \quad (i = n-1, 1) \\ R &= a_0 + Nb_1 \end{aligned}}$$

Newton's method for polynomials

18

- Newton's method works normally for simple (i.e. non repeated) roots.
- For polynomials it can be desirable to deflate the polynomial once a root is found

$$P_n(x) = (x - \alpha_1) Q_{n-1}(x)$$

- The remaining roots are then those of $Q_{n-1}(x)$
- The roots of $Q_2(x)$ can be found via quadratic formula. This offers a straight forward way to converge on all roots...
- Multiple roots are handled fine via Newton's method & there are some algorithms to speed convergence (see 3.S.2.3.)
 - a) multiplicity factor
 - b) modified function in Newton's method
- For multiple roots need to insure $f'(r) \neq 0 \Rightarrow$ use a relative convergence criteria based on $\underline{\Delta x = x_{i+1} - x_i}$
- For complex roots Newton's method can converge if an initial complex guess is chosen. (Bracketing methods do not work)
- Often deflated root solutions are refined through an additional iteration step.

Systems of nonlinear equations

(19)

- Newton's method can readily be adapted to solve multi-variate problems, e.g.:

$$f(x, y) = 0 \quad ; \quad g(x, y) = 0$$

- Expanding in a two-variable Taylor series:

$$f(x, y) = f(x_i, y_i) + \frac{\partial f}{\partial x}(x_i, y_i)(x - x_i)$$

$$+ \frac{\partial f}{\partial y}(x_i, y_i)(y - y_i) + \dots$$

$$g(x, y) = g_i + \left[\frac{\partial g}{\partial x} \right]_i (x - x_i) + \left[\frac{\partial g}{\partial y} \right]_i (y - y_i) + \dots$$

- To solve this system, we wish to drive f, g to zero
 $\Rightarrow f(x_i), g(x_i) \rightarrow 0$.

$$\begin{bmatrix} (\partial f / \partial x)_i & (\partial f / \partial y)_i \\ (\partial g / \partial x)_i & (\partial g / \partial y)_i \end{bmatrix} \begin{bmatrix} \Delta x \\ \Delta y \end{bmatrix} = \begin{bmatrix} -f_i \\ -g_i \end{bmatrix}$$

followed by the simple update $x_{i+1} = x_i + \Delta x$
 $y_{i+1} = y_i + \Delta y$

- More generally, for n equations, n unknowns:

(20)

$$f_j(x_k) = 0 \quad \begin{cases} \{x_k\} : \text{unknowns} \\ \{f_j\} : \text{nonlinear equations.} \end{cases}$$

- Define a matrix:

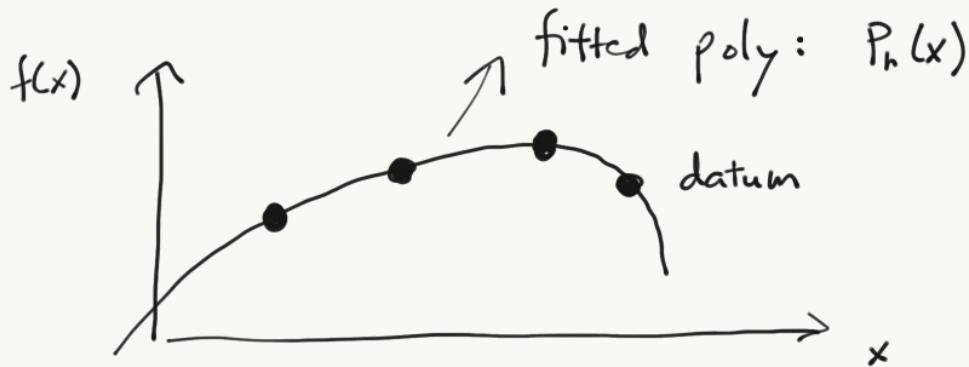
$$D_{jk} = \frac{\partial f_j}{\partial x_k}$$

$$\begin{aligned} \underline{\Delta x} &= -\underline{f} \\ \underline{x}_{ith} &= \underline{x}_i + \underline{\Delta x} \end{aligned} \quad \left. \right\} \text{So @ each iteration a linear system must be solved!}$$

Polynomials & fitting

①

- The most common applications polynomials is for fitting data or performing interpolation.
- They are easy to differentiate & integrate & inexpensive to compute
- Fitting can be either exact (passing thru all available data) or approximate (minimizing some error or residual metric)
- Most straightforward use of polynomials is direct fitting: viz. find the n^{th} degree polynomial exactly passing through $n+1$ points.



$$\underbrace{P_n(x)}_{f(x)} = a_0 + a_1 x + \dots + a_n x^n = \sum_{i=1}^n a_i x^i$$

$$f_j = \sum_{i=1}^n a_i x^i \rightarrow \underbrace{(\underline{x}^n)^T}_{\text{powers of } x} \underline{a} \quad (\text{matrix mult.})$$

↓
coeff vector

$$\underline{f} = \underline{X} \underline{a} \rightarrow \text{solve for } \underline{a} \quad \underline{X} = \begin{bmatrix} 1 & x_1 & x_1^2 & \dots \\ 1 & x_2 & x_2^2 & \dots \\ 1 & x_3 & x_3^2 & \dots \end{bmatrix}$$

Approximate fits

(9)

- In some situations we have effectively more data than unknowns such that the system is over-determined.

under-determined: fewer eqns. than unknowns
(need prior info to solve)

even-determined: same # of eqns. & unknowns

over-determined: more eqns. than unknowns.

- For an over determined problem we seek to fit the data in the best possible way. Letting y_i be the data & $f(x)$ over objective fn.

$$e_i = y_i - f(x_i) \quad (\text{error/deviation})$$

- least squares minimizes squared error: e_i^2

- linear objective fns. often used:

$$f(x) = a + bx ; \text{ need to minimize } S = \sum_i (y_i - f(x_i))^2$$

where S is understood to be a fn. of a, b .

$$1) \frac{\partial S}{\partial a} = \sum_i 2(y_i - f(x_i))(-1) \quad \left. \right\} \text{ set to zero}$$

$$2) \frac{\partial S}{\partial b} = \sum_i 2(y_i - f(x_i))(-x_i) \quad \left. \right\} \text{ for extrema}$$

$$1) \sum_i y_i = \sum_i a + bx_i \quad 2) \sum_i x_i y_i = \sum_i (ax_i + bx_i^2)$$

(To)

- Note that there are just two unknowns to solve for:

$$1) \quad a \sum_i 1 + b \sum_i x_i = \sum_i y_i$$

$$2) \quad a \sum_i x_i + b \sum_i x_i^2 = \sum_i x_i y_i$$

$$\underline{M} \underline{a} = \underline{y}'$$

$$\underline{M} = \begin{bmatrix} \sum_i 1 & \sum_i x_i \\ \sum_i x_i & \sum_i x_i^2 \end{bmatrix} \quad \underline{a} = \begin{bmatrix} a \\ b \end{bmatrix}$$

$$\underline{y}' = \begin{bmatrix} \sum_i y_i \\ \sum_i x_i y_i \end{bmatrix} \quad (\text{matrix system to be solved}) \dots$$

- Note that the original system of equations can be viewed itself in the context of matrix multiplications:

$$f(x_i) = a + b x_i \rightsquigarrow f_i = \underline{x}_i^T \underline{a}$$

$$\underline{x}_i^T = [1 \ x_i] \quad \underline{a}^T = [a \ b]$$

For the entire set of points: $\underline{f} = \underline{X} \underline{a}$

$$\underline{X} = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}; \quad \underline{f} = [f_i]; \quad \underline{a}^T = [a \ b]$$

(11)

- Note that, due to linearity of f we may write:

$$\underline{f} = \begin{bmatrix} \frac{\partial f}{\partial a} & \frac{\partial f}{\partial b} \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} \quad (\text{Jacobian matrix})$$

$$\underline{f} = \begin{bmatrix} \frac{\partial f}{\partial a}|_1 & \frac{\partial f}{\partial b}|_1 \\ \vdots & \vdots \\ \frac{\partial f}{\partial a}|_n & \frac{\partial f}{\partial b}|_n \end{bmatrix} \begin{bmatrix} a \\ b \end{bmatrix} = \underline{J} \underline{a}$$

- Because this is not a square system there is no direct solution for \underline{a}
- The error is:

$$\underline{e} = \{e_i\} = \underline{Y} - \underline{f} \quad (\text{vector of errors})$$

- The total squared error $\sum_i e_i^2$ is then:

$$\underline{e}^T \underline{e} = (\underline{Y} - \underline{J} \underline{a})^T (\underline{Y} - \underline{J} \underline{a}) = \sum_i e_i^2$$

- And the system of eqns. to be solved is:

$$\frac{\partial}{\partial a_j} (\underline{e}^T \underline{e}) = 0 \quad (j \text{ runs over unknowns } (1, 2))$$

i.e. # of data points j = # of unknowns.

$$\frac{\partial}{\partial a_j} (\underline{e}^T \underline{e}) = \frac{\partial}{\partial a_j} (\underline{e}^T) \underline{e} + \underline{e}^T \frac{\partial \underline{e}}{\partial a_j}$$

note: $\frac{\partial \underline{e}}{\partial a_j} = -\underline{J}_j$ (j th column of \underline{J} times a_j)

- Similarly :

$$\frac{\partial \underline{e}^T}{\partial \underline{a}_j} = \underline{J}_j^T \quad \left| \begin{array}{l} \underline{a} \sim mx1 \quad \underline{Y} \sim nx1 \\ \underline{e} \sim nx1 \quad \underline{J} \sim n \times m \end{array} \right.$$

$$\frac{\partial}{\partial \underline{a}_j} (\underline{e}^T \underline{e}) = \underline{J}_j^T \underline{e} + \underline{e}^T (\underline{J}_j) \quad \text{conform to add!!!}$$

- The full system to be solved is given by:

$$\leftarrow \frac{\partial}{\partial \underline{a}} (\underline{e}^T \underline{e}) = (\underline{J}^T) \underline{e} + (\underline{e}^T (\underline{J}))^T = 0$$

\rightarrow vector derivative, viz. gradient

$$\underline{J}^T (\underline{Y} - \underline{J} \underline{a}) + ((\underline{Y} - \underline{J} \underline{a})^T \underline{J})^T = 0$$

$$\underline{J}^T \underline{Y} - \underline{J}^T \underline{J} \underline{a} + \underline{J}^T (\underline{Y} - \underline{J} \underline{a}) = 0$$

$$\underline{J}^T \underline{Y} - \underline{J}^T \underline{J} \underline{a} + \underline{J}^T \underline{Y} - \underline{J}^T \underline{J} \underline{a} = 0$$

$$2(\underline{J}^T \underline{J}) \underline{a} = 2 \underline{J}^T \underline{Y}$$

$$(\underline{J}^T \underline{J}) \underline{a} = \underline{J}^T \underline{Y}$$

$$\therefore \underline{a} = (\underline{J}^T \underline{J})^{-1} \underbrace{\underline{J}^T \underline{Y}}$$

Moore-Penrose inverse

(solves over
determined systems)

Lagrange Polynomials

(2)

- Given two points $[a, f(a)] [b, f(b)]$ the line passing thru these is given by :

$$P_1(x) = \underbrace{\frac{(x-b)}{(a-b)} f(a)}_{\textcircled{O} \text{ for } x=b} + \underbrace{\frac{(x-a)}{(b-a)} f(b)}_{\textcircled{O} \text{ for } x=a}$$

- Note that :

$$P_1(a) = \frac{(a-b)}{(a-b)} f(a) = f(a) \quad \checkmark$$

$$P_1(b) = f(b) \quad \checkmark$$

- Consider three points & the quadratic form :

$$a, b, c \rightarrow f(a), f(b), f(c)$$

$$P_2(x) = \frac{(x-b)(x-c)}{(a-b)(a-c)} f(a) + \frac{(x-a)(x-c)}{(b-a)(b-c)} f(b) \\ + \frac{(x-a)(x-b)}{(c-a)(c-b)} f(c)$$

- This polynomial has similar properties : (1) coeff. of $f(a)$ is zero for $x = b$ or c , (2) coeff. of $f(a) = 1$ for $x = a$.

- Extension to higher order polynomial :

$$P_n(x) = \sum_{j=1}^{n+1} \frac{\prod_{i \neq j} (x - \alpha_i)}{\prod_{i \neq j} (\alpha_j - \alpha_i)} f(\alpha_j)$$

NOTE : poly representations are unique so that this gives same result as direct fit !

(3)

- Lagrange polynomials can be constructed sequentially thru use of differences:

$$f(x) = \frac{(x-b)}{(a-b)} f(a) + \frac{(x-a)}{(b-a)} f(b)$$

$$= \frac{f(b)(x-a) - f(a)(x-b)}{(b-a)}$$

- For two points $[a, b] \rightarrow [x_i, x_{i+1}]$

$$f_i^{(1)} = \frac{f(x_{i+1})(x - x_i) - f(x_i)(x - x_{i+1})}{(x_{i+1} - x_i)}$$

- For three points $[x_i, x_{i+1}, x_{i+2}]$:

$$f_i^{(2)} = \frac{(x - x_i) f_{i+1}^{(1)} - (x - x_{i+2}) f_i^{(1)}}{x_{i+2} - x_i}$$

$$f_i^{(n)} = \frac{(x - x_i) f_{i+1}^{(n-1)} - (x - x_{i+n}) f_i^{(n-1)}}{x_{i+n} - x_i}$$

generally, this can be most easily seen
by direct substitution

- Computing coefficients this way is known as
Neville's Algorithm

(4)

- Divided differences & associated polynomials are another way of conceptualizing Lagrange polynomials:

$$f_i^{(1)} = \frac{f_{i+1} - f_i}{x_{i+1} - x_i} \quad (\text{derivative approximation})$$

$$f_i^{(2)} = \frac{f_{i+1}^{(1)} - f_i^{(1)}}{x_{i+2} - x_i} \quad (\text{2nd derivative})$$

$$= \frac{\frac{f_{i+2} - f_{i+1}}{x_{i+2} - x_{i+1}} - \frac{f_{i+1} - f_i}{x_{i+1} - x_i}}{x_{i+2} - x_i}$$

$$= f_{i+2} \left(\frac{1}{(x_{i+2} - x_{i+1})(x_{i+2} - x_i)} \right)$$

$$+ f_{i+1} \left(\frac{-1}{(x_{i+2} - x_{i+1})(x_{i+2} - x_i)} - \frac{1}{(x_{i+1} - x_i)(x_{i+2} - x_i)} \right)$$

$$+ f_i \left(\frac{1}{(x_{i+1} - x_i)(x_{i+2} - x_i)} \right)$$

- These can also be used to construct polynomials that pass thru $n+1$ pts. \rightarrow (of n^{th} degree)

$$f(x) = f_i^{(0)} + f_i^{(1)}(x - x_i) + f_i^{(2)}(x - x_{i+1})(x - x_i) + \dots + f_i^{(n)}(x - x_{i+n-1})(x - x_{i+n-2}) \dots (x - x_i)$$

(Finite) Differences

(5)

- Forward difference : $\Delta f_i = (f_{i+1} - f_i)$
- Backward diff : $\bar{\Delta} f_i = (f_i - f_{i-1})$
- Centered diff : $\delta f_i = (f_{i+1/2} - f_{i-1/2})$
OR $\delta f_{i+1/2} = (f_{i+1} - f_i)$
- The book makes frequent use of difference tables:

| <u>x</u> | <u>f(x)</u> | <u>Δf</u> | <u>$\Delta^2 f$</u> | <u>$\Delta^3 f$</u> |
|----------|-------------|------------------------------|--------------------------------|--------------------------------|
| x_0 | f_0 | --- | Δf_0 | |
| x_1 | f_1 | Δf_0 | --- | $\Delta^2 f_0$ |
| x_2 | f_2 | Δf_1 | $\Delta^2 f_0$ | --- |
| x_3 | f_3 | Δf_2 | $\Delta^2 f_1$ | $\Delta^3 f_0$ |

- Higher order differences are defined by :

$$\Delta^n f_i = \Delta^{(n-1)} f_{i+1} - \Delta^{(n-1)} f_i$$

so that $\Delta^2 f_i = \Delta f_{i+1} - \Delta f_i$

$$= f_{i+2} - f_{i+1} - (f_{i+1} - f_i)$$

$$= f_{i+2} - 2f_{i+1} + f_i$$

and so on ...

- Note that errors get compounded as one goes across the table and computes higher order differences.
- The Newton forward difference polynomial is defined by :

$$\begin{aligned}
 P_h(x) = f_0 + s \Delta f_0 + \frac{s(s-1)}{2!} \Delta^2 f_0 \\
 + \frac{s(s-1)(s-2)}{3!} \Delta^3 f_0 + \dots \\
 + \frac{s(s-1)(s-2) \dots (s-(n-1))}{n!} \Delta^n f_0
 \end{aligned}$$

where $s = \frac{x - x_0}{\Delta x} = \frac{x - x_0}{h} \Rightarrow x = x_0 + sh$

$(s \in \mathbb{Z}^+)$ (set of positive integers)

- Coefficients of this polynomial are binomial coefficients:

$$\frac{s(s-1)(s-2) \dots (s-(n-1))}{n!} = \frac{s!}{n! (s-n)(s-n-1) \dots (1)}$$

$$= \left[\frac{s!}{n!(s-n)!} = \binom{s}{n} \right] \quad \text{"choose n"}$$

- Other forms of difference polynomials also exist (see book).

Multivariate Approximation

(7)

- We are often faced with fitting or approximation of functions of more than one variable.

e.g. $\underline{z = f(x,y)}$

usually a multivariate polynomial is a good choice:

$$f(x,y) = a + bx + cy + dx\bar{y} + ex^2 + fy^2 \\ + gx^2y + hy^2x + \dots$$

- A direct fit approach can work for lower-degree functions $f(x,y)$:
 $f(x,y) \approx a + bx + cy + d\bar{xy}$ (bilinear interpolant)
- For a given set of four data points we can produce a solution for $\{a, b, c, d\}$

$$\begin{bmatrix} f_1 \\ f_2 \\ f_3 \\ f_4 \end{bmatrix} = \begin{bmatrix} 1 & x_1 & y_1 & x_1 y_1 \\ 1 & x_2 & y_2 & x_2 y_2 \\ 1 & x_3 & y_3 & x_3 y_3 \\ 1 & x_4 & y_4 & x_4 y_4 \end{bmatrix} \begin{bmatrix} a \\ b \\ c \\ d \end{bmatrix}$$

which can be solved thru standard linear algebraic approaches.

- Solution for higher-order coefficients will be ill-⁽⁸⁾ conditioned.
- Successive univariate approximation is often a better approach:

$$f_{j*}(x) = f(x, y_j) \approx a_j + b_j x + c_j x^2 + d_j x^3 + \dots$$

(one of these
for each point)

denote: $f_{1j} \Rightarrow f(x_1, y_j)$

$$\begin{bmatrix} f_{1j} \\ f_{2j} \\ f_{3j} \\ \vdots \\ f_{nj} \end{bmatrix} = \begin{bmatrix} 1 & x_1 & \dots & x_1^{n-1} \\ 1 & x_2 & & x_2^{n-1} \\ 1 & x_3 & & x_3^{n-1} \\ \vdots & \vdots & & \vdots \\ 1 & x_n & \dots & x_n^{n-1} \end{bmatrix} \begin{bmatrix} a_i \\ b_i \\ c_i \\ \vdots \end{bmatrix}$$

- which can be solved for the coefficients a_j, b_j etc
- Once we have interpolating fns. for $f(x, y_j)$ we can fit these to polynomials in the y variable:

$$f(x_i, y) = a'_i + b'_i y + c'_i y^2 + \dots$$

evaluated using $f_{j*}(x)$, then set up a system to solve:

approximating f_n .
is given from \underline{a}'

$$\left\{ \begin{array}{l} f_* = \underline{\underline{Y}} \underline{a}' \rightarrow \boxed{\underline{a}' = \underline{\underline{Y}}^{-1} \cdot \underline{\underline{f}_x}} \\ (\text{probably use Gauss elim}) \end{array} \right.$$

Splines

(13)

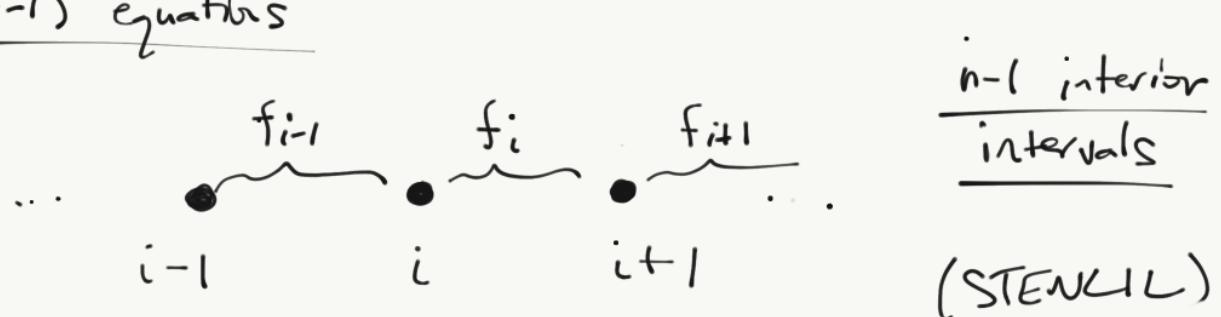
- Splines are effectively piecewise polynomials fitted to a subset of available data points.
- Cubics are the most commonly used as their slope & curvature can be forced to be continuous @ each data point.

$$f_i(x) = a_i + b_i x + c_i x^2 + d_i x^3 \quad (i=1, n+1)$$

- Cubic Splines satisfy the constraints:

↓
(total #
of data
pts.)

- 0) $f_i(x)$ must pass through the data point @ $x=x_i$: $n-1$ constraints. (interior)
- 1) $f_i(x)$ & $f_{i+1}(x)$ have the same value @ $x=x_i$: $(n-1)$ constraints.
- 2) 1st derivatives are equal on each side of a given pt. (interior) : $f'_i(x) = f'_{i+1}(x)$ @ $x=x_i$
 $(n-1)$ equations
- 3) 2nd derivatives of adjacent splines are equal @ interior grid points : $f''_i(x) = f''_{i+1}(x)$ @ $x=x_i$
 $(n-1)$ equations



- There remains four additional constraints to specify
a unique solution. (14)



- 4) f_i must pass thru the data @ x_i & f_n must pass thru the data @ x_{n+1} : 2 equations
- 5) curvature is specified @ the first & last grid points: 2 equations

Numerical Differentiation

①

- One approach to differentiate a function is to perform differentiation on an approximation (like a polynomial):

$$\frac{dF}{dx} = \frac{d}{dx}(P_n(x))$$

- Integration can similarly be performed in this manner

$$\int f(x) dx = \int P_n(x) dx$$

- For differentiation it is often most enlightening to use a Taylor Series approach:

$$f(x) = \sum_n f^{(n)}(x_0) \frac{(x-x_0)^n}{n!}$$

$$= f(x_0) + f'(x_0)(x-x_0) + \frac{f''(x_0)}{2!}(x-x_0)^2 + \dots$$

- Letting $x = x_0 + \Delta x$

$$\textcircled{a} \quad f(x+\Delta x) = f(x_0) + f'(x_0) \Delta x + \frac{f''(x_0)}{2!} \Delta x^2 + \frac{f'''(x_0)}{3!} \Delta x^3$$

- Similarly expanding for $x - \Delta x$:

$$\textcircled{b} \quad f(x-\Delta x) = f(x_0) - f'(x_0) \Delta x + \frac{f''(x_0)}{2!} \Delta x^2 - \frac{f'''(x_0)}{3!} \Delta x^3 + \dots$$

- These may be combined (a - b):

$$f(x+\Delta x) - f(x-\Delta x) = 2 f'(x_0) \Delta x + \frac{2}{3!} f'''(x_0) \Delta x^3$$

$$f'(x_0) = \frac{f(x_0 + \Delta x) - f(x_0 - \Delta x)}{2\Delta x} + \frac{1}{6} f'''(x_0) \Delta x^2$$

= " + $\mathcal{O}(\Delta x^2)$

error term varying
 like Δx^2

- These two eqns. (a, b) can be added to yield an approx. 2nd derivative:

$$f(x_0 + \Delta x) + f(x_0 - \Delta x) = 2f(x_0) + f''(x_0) \Delta x^2 + \frac{1}{4!} f'''(x_0) \Delta x^4 + \dots$$

$$\therefore f''(x_0) \approx \frac{f(x_0 + \Delta x) - 2f(x_0) + f(x_0 - \Delta x)}{\Delta x^2} + \mathcal{O}(\Delta x^2)$$

- Often use indices to represent locations on a grid

$$x_0 \rightarrow i ; \quad x_0 + \Delta x \rightarrow i+1 ; \quad x_0 - \Delta x \rightarrow i-1$$

$$\left[\frac{df}{dx} \right]_i = \frac{f_{i+1} - f_{i-1}}{2\Delta x}$$

$$\left[\frac{d^2f}{dx^2} \right]_i = \frac{f_{i+1} - 2f_i + f_{i-1}}{\Delta x^2}$$

for n^{th} order derivative generally need to solve n equations from Taylor Series

$\cdots - i-1 - i \bullet - i+1 - \cdots$

\rightarrow (Stencil of FD approx.)

- Less accurate approximations are possible :

(3)

$$f(x) = f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x-x_0)^2}{2!} + \dots$$

$$f(x+\Delta x) = f(x_0) + f'(x_0)\Delta x + \frac{f''(\Delta x)^2}{2!} + \dots$$

$$f'(x_0) = \frac{f(x_0 + \Delta x) - f(x_0)}{\Delta x} + \frac{f''(\Delta x)}{2!} + \dots$$

(forward difference f' approx..)

- Also can use a different expansion to estimate f' :

$$f(x_0 - \Delta x) = f(x_0) - f'(x_0)\Delta x + \frac{f''(x_0)\Delta x^2}{2!} + \dots$$

$$\frac{f(x_0) - f(x_0 - \Delta x)}{\Delta x} = f'(x_0) + \frac{f''(\Delta x)}{2!} + \dots$$

(backward difference)

- The Taylor Series approach is generally applicable to produce difference formulas for derivatives of any order at any accuracy level.
- For example we may develop an expression for the third derivative:

$$f(x+\Delta x) = f(x) + f'(x)\Delta x + \frac{1}{2!} f''(x)\Delta x^2 + \frac{1}{3!} f'''(x)\Delta x^3 + \frac{1}{4!} f''''(x)\Delta x^4 + \dots$$

- For a third derivative estimate we need to remove f', f'' from these equations. (it will have $\mathcal{O}(\Delta x)$ accuracy) (4)

$$(1) \quad f_{i+2} = f_i + f' 2\Delta x + \frac{f''}{2!} 4\Delta x^2 + \frac{f'''}{3!} 8\Delta x^3 + \mathcal{O}(\Delta x^4)$$

$$(2) \quad f_{i+1} = f_i + f' \Delta x + \frac{f''}{2!} \Delta x^2 + \frac{f'''}{3!} \Delta x^3 + \mathcal{O}(\Delta x^4)$$

$$(3) \quad f_{i-1} = f_i - f' \Delta x + \frac{f''}{2!} \Delta x^2 - \frac{f'''}{3!} \Delta x^3 + \mathcal{O}(\Delta x^4)$$

- We aim to eliminate f', f'' from these to solve for f''' (so we have 3 eqns. & 3 unknowns). $f_{i+2}, f_{i+1}, f_i, f_{i-1}$ are assumed to be given:

- Alternatively we could eliminate $\mathcal{O}(\Delta x^4)$ by adding another eqn.

$$\mathcal{O}(\Delta x^4) = \frac{f^{(4)}}{4!} \Delta x^4 + \frac{f^{(5)}}{5!} \Delta x^5 + \dots$$

(in eqn for f_{i+1})

$$(4) \quad f_{i-2} = f_i - 2f' \Delta x + \frac{f''}{2!} 4\Delta x^2 - \frac{f'''}{3!} 8\Delta x^3 + \frac{f^{(4)}}{4!} 16\Delta x^4 + \mathcal{O}(\Delta x^5)$$

- Now you have 4 eqns. & 4 unknowns.
↳ $(f', f'', f''', f^{(4)})$

(5)

- Finite difference formulas can be applied iteratively to derive higher-order differences. E.g.:

$$f' \approx \frac{f_{i+1} - f_{i-1}}{2\Delta x} + O(\Delta x^2)$$

$$f'' \approx \frac{(f')_{i+1/2} - (f')_{i-1/2}}{\Delta x} \quad \left. \begin{array}{l} \text{some approximation} \\ \text{of derivative @} \\ \text{half point (2nd} \\ \text{order accurate)} \end{array} \right\}$$

$$f'_{i+1/2} = \frac{f_{i+1} - f_i}{\Delta x} \quad (\text{centered formula})$$

$$f'' \approx \frac{\frac{f_{i+1} - f_i}{\Delta x} - \frac{f_i - f_{i-1}}{\Delta x}}{\Delta x}$$

- This approach allows differentiation of complicated expressions like:

$$\frac{d}{dx} \left(A(x) \frac{df}{dx} \right) \quad \& \text{ similar.}$$

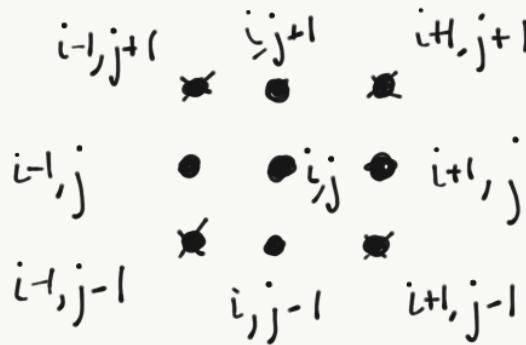
- Iterative applications of difference formula are also useful for deriving mixed derivatives

$$\left[\frac{\partial^2 f}{\partial x \partial y} \right]_{i,j} \approx \frac{1}{2\Delta x} \left(\left[\frac{\partial f}{\partial y} \right]_{i+1,j} - \left[\frac{\partial f}{\partial y} \right]_{i-1,j} \right)$$

⑥

$$\left[\frac{\partial f}{\partial y} \right]_{i+1} \approx \frac{(f_{i+1, j+1} - f_{i+1, j-1})}{2 \Delta y}$$

$$\therefore \left[\frac{\partial^2 f}{\partial x \partial y} \right]_{i,j} = \frac{1}{4 \Delta x \Delta y} (f_{i+1, j+1} - f_{i+1, j-1} - f_{i-1, j+1} + f_{i-1, j-1})$$



other approximations
are possible this one
preserves:

$$\left[\frac{\partial^2 f}{\partial x \partial y} \right]_{i,j} = \left[\frac{\partial^2 f}{\partial y \partial x} \right]_{i,j}$$

- Another approach would be to just differentiate iteratively :